

**CLAIMS:**

What is claimed is:

1 1. A method of selecting data sets for use with a  
2 predictive algorithm based on data network geographical  
3 information, comprising:

4 generating a first distribution of a training data  
5 set;

6 generating a second distribution of a testing data  
7 set;

8 comparing the first distribution and the second  
9 distribution to identify a discrepancy between the first  
10 distribution and the second distribution with respect to  
11 data network geographical information; and

12 modifying selection of entries in one or more of the  
13 training data set and the testing data set based on the  
14 discrepancy between the first distribution and the second  
15 distribution.

1 2. The method of claim 1, wherein the first  
2 distribution and the second distribution are  
3 distributions of a number of data network links from a  
4 customer data network geographical location to a web site  
5 data network geographical location.

1 3. The method of claim 1, wherein the first  
2 distribution and the second distribution are  
3 distributions of a size of a click stream for arriving at  
4 a web site data network geographical location.

1 4. The method of claim 1, wherein comparing the first  
2 distribution and the second distribution includes  
3 comparing one or more of a mean, mode, and standard  
4 deviation of the first distribution to one or more of a  
5 mean, mode, and standard deviation of the second  
6 distribution.

1 5. The method of claim 1, wherein the first  
2 distribution and the second distribution are  
3 distributions of a weighted data network geographical  
4 distance between a customer data network geographical  
5 location and a web site data network geographical  
6 locations.

1 6. The method of claim 1, wherein the first  
2 distribution and the second distribution are  
3 distributions of a weighted click stream for arriving at  
4 a web site data network geographical locations.

00000000000000000000000000000000

1 7. The method of claim 1, wherein modifying selection  
2 of entries in one or more of the training data set and  
3 the testing data set includes generating recommendations  
4 for improving selection of entries in one or more of the  
5 training data set and the testing data set.

1 8. The method of claim 1, wherein the training data set  
2 and the testing data set are selected from a customer  
3 information database.

1 9. The method of claim 1, further comprising comparing  
2 at least one of the first distribution and the second  
3 distribution to a distribution of a customer database.

1 10. The method of claim 1, wherein the first  
2 distribution and second distribution are frequency  
3 distributions of one of number of data network links  
4 between a customer geographical location and one or more  
5 web site data network geographical locations, and size of  
6 a click stream for arriving at one or more web site data  
7 network geographical locations.

1 11. The method of claim 9, wherein comparing at least  
2 one of the first distribution and the second distribution

3 to a distribution of a customer database includes:

4 generating a composite data set from the training  
5 data set and the testing data set; and

6 generating a composite distribution from the  
7 composite data set.

1 12. The method of claim 1, wherein modifying selection  
2 of entries in one or more of the training data set and  
3 the testing data set includes changing one of a random  
4 selection algorithm and a seed value for a random  
5 selection algorithm.

1 13. The method of claim 1, further comprising training a  
2 predictive algorithm using at least one of the training  
3 data set and the testing data set if the discrepancy is  
4 within a predetermined tolerance.

1 14. The method of claim 13, wherein the predictive  
2 algorithm is a discovery based data mining algorithm.

1 15. An apparatus for selecting data sets for use with a  
2 predictive algorithm based on data network geographical  
3 information, comprising:

4 a statistical engine; and

5 a comparison engine coupled to the statistical  
6 engine, wherein the statistical engine generates a first  
7 distribution of a training data set and a second  
8 distribution of a testing data set, the comparison engine  
9 compares the first distribution and the second  
10 distribution to identify a discrepancy between the first  
11 distribution and the second distribution with respect to  
12 data network geographical information, and modifies  
13 selection of entries in one or more of the training data  
14 set and the testing data set based on the discrepancy  
15 between the first distribution and the second  
16 distribution.

1 16. The apparatus of claim 15, wherein the first  
2 distribution and the second distribution are  
3 distributions of a number of data network links from a  
4 customer data network geographical location to a web site  
5 data network geographical location.

1 17. The apparatus of claim 15, wherein the first  
2 distribution and the second distribution are  
3 distributions of a size of a click stream to arrive at a  
4 web site data network geographical location.

1234567890

1 18. The apparatus of claim 15, wherein the comparison  
2 engine compares the first distribution and the second  
3 distribution by comparing one or more of a mean, mode,  
4 and standard deviation of the first distribution to one  
5 or more of a mean, mode, and standard deviation of the  
6 second distribution.

1 19. The apparatus of claim 15, wherein the first  
2 distribution and the second distribution are  
3 distributions of a weighted number of data network links  
4 between a customer data network geographical location and  
5 a web site data network geographical location.

1 20. The apparatus of claim 15, wherein the first  
2 distribution and the second distribution are  
3 distributions of a weighted size of a click stream to  
4 arrive at a web site data network geographical location.

1 21. The apparatus of claim 15, wherein the comparison  
2 engine modifies selection of entries in one or more of  
3 the training data set and the testing data set by  
4 generating recommendations for improving selection of  
5 entries in one or more of the training data set and the  
6 testing data set.

1 22. The apparatus of claim 15, further comprising a  
2 training data set/testing data set selection device that  
3 selects the training data set and the testing data set  
4 from a customer information database.

1 23. The apparatus of claim 15, wherein the comparison  
2 engine further compares at least one of the first  
3 distribution and the second distribution to a  
4 distribution of a customer database.

1 24. The apparatus of claim 15, wherein the first  
2 distribution and second distribution are frequency  
3 distributions of one of a number of data network links  
4 between a customer data network geographical location and  
5 one or more web site data network geographical locations,  
6 and a size of a click stream to arrive at one or more web  
7 site data network geographical locations.

1 25. The apparatus of claim 23, wherein the comparison  
2 engine compares at least one of the first distribution  
3 and the second distribution to a distribution of a  
4 customer database by:

1 23. The apparatus of claim 15, wherein the first  
2 distribution and second distribution are frequency  
3 distributions of one of a number of data network links  
4 between a customer data network geographical location and  
5 one or more web site data network geographical locations,  
6 and a size of a click stream to arrive at one or more web  
7 site data network geographical locations.

5 generating a composite data set from the training  
6 data set and the testing data set; and

7 generating a composite distribution from the  
8 composite data set.

1 26. The apparatus of claim 15, wherein the comparison  
2 engine modifies selection of entries in one or more of  
3 the training data set and the testing data set by  
4 changing one of a random selection algorithm and a seed  
5 value for a random selection algorithm.

1 27. The apparatus of claim 15, further comprising a  
2 predictive algorithm device, wherein the predictive  
3 algorithm device is trained using at least one of the  
4 training data set and the testing data set if the  
5 discrepancy is within a predetermined tolerance.

1 28. The apparatus of claim 27, wherein the predictive  
2 algorithm is a discovery based data mining algorithm.

1 29. A computer program product in a computer readable  
2 medium for selecting data sets for use with a predictive  
3 algorithm based on data network geographical information,  
4 comprising:

5        first instructions for generating a first  
6 distribution of a training data set;

7        second instructions for generating a second  
8 distribution of a testing data set;

9        third instructions for comparing the first  
10 distribution and the second distribution to identify a  
11 discrepancy between the first distribution and the second  
12 distribution with respect to data network geographical  
13 information; and

14       fourth instructions for modifying selection of  
15 entries in one or more of the training data set and the  
16 testing data set based on the discrepancy between the  
17 first distribution and the second distribution.

1 30. The computer program product of claim 29, wherein  
2 the first distribution and the second distribution are  
3 distributions of a number of data network links from a  
4 customer data network geographical location to a web site  
5 data network geographical location.

1 31. The computer program product of claim 29, wherein  
2 the first distribution and the second distribution are  
3 distributions of a size of a click stream to arrive at a  
4 web site data network geographical location.

1 32. The computer program product of claim 29, wherein  
2 the third instructions for comparing the first  
3 distribution and the second distribution include  
4 instructions for comparing one or more of a mean, mode,  
5 and standard deviation of the first distribution to one  
6 or more of a mean, mode, and standard deviation of the  
7 second distribution.

1 33. The computer program product of claim 29, wherein  
2 the first distribution and the second distribution are  
3 distributions of a weighted number of data network links  
4 between a customer data network geographical location and  
5 a web site data network geographical location.

1 34. The computer program product of claim 29, wherein  
2 the first distribution and the second distribution are  
3 distributions of a weighted size of a click stream to  
4 arrive at a web site data network geographical location.

1 35. The computer program product of claim 29, wherein  
2 the fourth instructions for modifying selection of  
3 entries in one or more of the training data set and the  
4 testing data set include instructions for generating  
5 recommendations for improving selection of entries in one  
6 or more of the training data set and the testing data  
7 set.

0002794341-050001

1 36. The computer program product of claim 29, further  
2 comprising fifth instructions for comparing at least one  
3 of the first distribution and the second distribution to  
4 a distribution of a customer database.

1 37. The computer program product of claim 29, wherein  
2 the first distribution and second distribution are  
3 frequency distributions of one of a number of data  
4 network links between a customer data network  
5 geographical location and one or more web site data  
6 network geographical locations, and a size of a click  
7 stream to arrive at one or more web site data network  
8 geographical locations.

1 38. The method of claim 36, wherein the fifth  
2 instructions include:  
3       instructions for generating a composite data set  
4       from the training data set and the testing data set; and  
5       instructions for generating a composite distribution  
6       from the composite data set.

09876543210987654

1 39. The computer program product of claim 29, wherein  
2 the fourth instructions for modifying selection of  
3 entries in one or more of the training data set and the  
4 testing data set include instructions for changing one of  
5 a random selection algorithm and a seed value for a  
6 random selection algorithm.

1 40. The computer program product of claim 29, further  
2 comprising fifth instructions for training a predictive  
3 algorithm using at least one of the training data set and  
4 the testing data set if the discrepancy is within a  
5 predetermined tolerance.

1 41. A method of predicting customer behavior based on  
2 data network geographical influences, comprising:  
3       obtaining data network geographical information  
4 regarding a plurality of customers;  
5       training a predictive algorithm using the data  
6 network geographical information; and  
7       using the predictive algorithm to predict customer  
8 behavior based on the data network geographical  
9 information.

1 42. An apparatus for predicting customer behavior based  
2 on data network geographical influences, comprising:

3 means for obtaining data network geographical  
4 information regarding a plurality of customers;

5 means for training a predictive algorithm using the  
6 data network geographical information; and

means for using the predictive algorithm to predict customer behavior based on the data network geographical information.

1 43. A computer program product in a computer readable  
2 medium for predicting customer behavior based on data  
3 network geographical influences, comprising:

4 first instructions for obtaining data network  
5 geographical information regarding a plurality of  
6 customers;

7 second instructions for training a predictive  
8 algorithm using the data network geographical  
9 information; and

third instructions for using the predictive algorithm to predict customer behavior based on the data network geographical information.